



POLITECNICO
MILANO 1863

SCUOLA DI INGEGNERIA INDUSTRIALE
E DELL'INFORMAZIONE

SMBUD Project - Spark

Author(s): **Gabriele Ginestroni**

Giacomo Gumiero

Lorenzo Iovine

Nicola Landini

Francesco Leone

Group Number: **10**

Academic Year: 2022-2023

Contents

Contents	i
1 Introduction	1
2 Data Structure	3
2.1 Article structure	3
2.2 Author structure	5
3 Dataframe Structure	7
3.1 Article Structure	7
3.2 Author Structure	8
3.3 Venue Structure	9
4 Commands and Queries	11
4.1 Commands	11
4.1.1 Insert a new author	11
4.1.2 Insert a new author	12
4.1.3 Insert a new venue	13
4.1.4 Insert a new article in his author dataframe	13
4.1.5 Update the number of citations of referenced publications	13
4.2 Queries	14
4.2.1 Query 1 - WHERE, JOIN	14
4.2.2 Query 2 - WHERE, LIMIT, LIKE	14
4.2.3 Query 3 - WHERE, IN, NESTED QUERY	15
4.2.4 Query 4 - GROUP BY, JOIN, AS	15
4.2.5 Query 5 - WHERE, GROUP BY	16
4.2.6 Query 6 - GROUP BY, HAVING, AS	16
4.2.7 Query 7 - WHERE, GROUP BY, HAVING, AS	16
4.2.8 Query 10 - WHERE, GROUP BY, HAVING, 2 JOINS	16

1 | Introduction

In this chapter will be presented the problem specification and the hypothesis under which the database is implemented.

2 | Data Structure

2.1. Article structure

In the following JSON file we can see how the *Article* is structured in our dataset.

Note: in some fields we decided not to insert all the content only for ease of read reasons.

```

1 {
2   "_id": "53e99f86b7602d9702859fdf",
3   "title": "Locality Sensitive Outlier Detection: A ranking driven
      approach",
4   "authors": [
5     {
6       "idAuth": "542a4c9fdabfae61d496694e",
7       "org": "Computer Science and Engineering Department, The Ohio State
      University, USA"
8     },
9     {
10      "idAuth": "53f48bc5dabfaea7cd1cce1d",
11      "org": "Computer Science and Engineering Department, The Ohio State
      University, USA"
12     },
13     {
14      "idAuth": "53f44b6fdabfaec09f1dd00d",
15      "org": "Computer Science and Engineering Department, The Ohio State
      University, USA"
16     }
17   ],
18   "n_citation": 60,
19   "abstract": "Outlier detection is fundamental to a variety of database
      and analytic tasks. Recently, distance-based outlier detection has
      emerged as a viable and scalable alternative to traditional
      statistical and geometric approaches. In this article we explore the
      role of ranking for the efficient discovery of distance-based
      outliers from large high dimensional data sets. Specifically, we
      develop a light-weight ranking scheme that is powered by locality

```

sensitive hashing, which reorders the database points according to their likelihood of being an outlier. We provide theoretical arguments to justify the rationale for the approach and subsequently conduct an extensive empirical study highlighting the effectiveness of our approach over extant solutions. We show that our ranking scheme improves the efficiency of the distance-based outlier discovery process by up to 5-fold. Furthermore, we find that using our approach the top outliers can often be isolated very quickly, typically by scanning less than 3% of the data set.",

```
20 "doi": "10.1109/ICDE.2011.5767852",
21 "keywords": [
22     "database point",
23     "ranking scheme",
24     "geometric approach",
25     ...
26 ],
27 "isbn": "978-1-4244-8958-9",
28 "page_start": "410",
29 "page_end": "421",
30 "year": 2011,
31 "fos": [
32     "Locality-sensitive hashing",
33     "Anomaly detection",
34     "Data mining",
35     ...
36 ],
37 "venue": {
38     "raw": "ICDE",
39     "type": 0
40 },
41 "references": [
42     "53e99fddb7602d97028b7e65"
43 ]
44 }
```


2.2. Author structure

In the following JSON file we can see how the *Author* is structured in our dataset.

Note: in some fields we decided not to insert all the content only for ease of read reasons.

```
1 {  
2   "_id": "53f45775dabfaee4dc8162e6",  
3   "name": "Guillermo Jorge-Botana",  
4   "articles": ["53e99f86b7602d970285a187"],  
5   "orcid": "0000-0001-5879-6783",  
6   "bio": "Qing-Long Han received the B.Sc. degree in mathematics from the  
       ...",  
7   "email": "Guillermo.Jorge-Botana@yahoo.com",  
8   "nationality": "de",  
9   "dob": "1945-06-19T00:00:00Z"  
10 }
```


3 | Dataframe Structure

3.1. Article Structure

```

1 root
2 |-- _id: string (nullable = true)
3 |-- title: string (nullable = true)
4 |-- authors: array (nullable = true)
5 |   |-- element: struct (containsNull = true)
6 |   |   |-- idAuth: string (nullable = true)
7 |   |   |-- org: string (nullable = true)
8 |-- n_citation: integer (nullable = true)
9 |-- abstract: string (nullable = true)
10 |-- doi: string (nullable = true)
11 |-- keywords: array (nullable = true)
12 |   |-- element: string (containsNull = true)
13 |-- isbn: string (nullable = true)
14 |-- page_start: string (nullable = true)
15 |-- page_end: string (nullable = true)
16 |-- year: integer (nullable = true)
17 |-- fos: array (nullable = true)
18 |   |-- element: string (containsNull = true)
19 |-- references: array (nullable = true)
20 |   |-- element: string (containsNull = true)
21 |-- issue: string (nullable = true)
22 |-- volume: string (nullable = true)
23 |-- publisher: string (nullable = true)
24 |-- venue_raw: string (nullable = true)

```

The structure just shown represents an *Article*; its attributes are:

- **_id** is the identifier of a publication.
- **title** represents the title of the publication.
- **authors** is an array that contains: **idAuth** of the authors of the article and the **org** field which represent the affiliation.

- **n_citation** is the number of times that the publication has been mentioned.
- **abstract** is a string containing a brief summary of the contents of the paper.
- **doi** Digital Object Identifier is a persistent and standardized identifier.
- **keywords** is an array containing keywords of the publication.
- **isbn** is an identification code of the venue of the publication.
- **page_start** defines the starting page of the publication.
- **page_end** defines the last page of the publication.
- **year** represents the year of publication.
- **fos** is an array containing the fields of study of the publication.
- **references** set of ObjectIds of the referenced articles.
- **issue** refers to how many times a periodical has been published during that year.
- **volume** is the volume of the venue in which the article has been published.
- **publisher** is the name of the publisher of the article.
- **venue_raw** is the name or the abbreviation of the venue (regardless the year, issue or volume) in which the publication was presented.

3.2. Author Structure

```

1 root
2 |-- _id: string (nullable = true)
3 |-- name: string (nullable = true)
4 |-- nationality: string (nullable = true)
5 |-- articles: array (nullable = true)
6 |   |-- element: string (containsNull = true)
7 |-- bio: string (nullable = true)
8 |-- email: string (nullable = true)
9 |-- orcid: string (nullable = true)
10 |-- dob: timestamp (nullable = true)

```

The structure just shown represents an *Author*; its attributes are:

- **_id** is the identifier of an author.
- **name** is the name of the author.

- **nationality** is the nationality of the author.
- **articles** is a set of articles identifier of the publications of the author.
- **bio** is a string that describes the author.
- **email** is the email address of the author.
- **orcid** Open Researcher and Contributor ID is a unique identifier for authors of scientific articles.
- **dob** is the birth date of the author.

3.3. Venue Structure

```
1 root
2 |-- raw: string (nullable = true)
3 |-- type: integer (nullable = true)
4 |-- artIds: array (nullable = false)
5 |     |-- element: string (containsNull = false)
6 |-- city: string (nullable = true)
```

The structure just shown represents a *Venue*. This dataframe was obtained using data imported from the article JSON shown before. *Venue* attributes are the following:

- **raw** is the name or the abbreviation of the venue (regardless the year, issue or volume) in which the publication was presented.
- **type** indicates the type of the publication.
- **artIds** is a set of articles identifier associated to the venue.
- **city** represents the location of the venue.

4.1.2. Insert a new author

Assuming it is not present in the dataset, we created a new row with the values for a new article written by the author created in section 4.1.1. In order to set the authors, we instantiated an array `new_authors`.

```

1 new_authors = [Row("638db170ae9ea0d19fad7a79", "Politecnico di Milano")
2               , Row("638db170ae9ea0d19fad7a7a", "
3                   Politecnico di Milano")]
4
5 new_article = Row(
6     _id="638db237d794b76f45c77916",
7     title="An extensive study of C-SMOTE, a Continuous Synthetic
8           Minority Oversampling Technique for
9           Evolving Data Streams",
10    authors=new_authors,
11    n_citation=3,
12    abstract = "Streaming Machine Learning (SML) studies algorithms that
13               update their models,\
14               given an unbounded and often non-stationary flow of data
15               performing a single pass. Online \
16               class imbalance learning is a branch of SML that combines the
17               challenges of both class imbalance\
18               and concept drift. In this paper, we investigate the binary
19               classification problem by
20               rebalancing\
21               an imbalanced stream of data in the presence of concept drift,
22               accessing one sample at a time.",
23    doi="10.1016/j.eswa.2022.116630",
24    keywords=["Evolving Data Stream","Streaming","Concept drift","
25              Balancing"],
26    isbn="123-4-567-89012-3",
27    page_start="39",
28    page_end="46",
29    year=2022,
30    fos=["Computer Science","Stream Reasoning","Big Data"],
31    references=["53e99fe4b7602d97028bf743","53e99fddb7602d97028bc085"],
32    issue="1",
33    volume="196",
34    publisher="Elsevier",
35    venue_raw="ESA"
36 )

```



```
27 df_articles = df_articles.union(spark.createDataFrame([new_article]))
```

4.1.3. Insert a new venue

Assuming it is not present in the dataset, we created a new row with the values for a new venue *ESA* hosted in *Montreal*.

Note: in field `artIds` we set the article created in section 4.1.2.

```
1 new_venue = Row(
2     raw="ESA",
3     type=1,
4     artIds=["638db237d794b76f45c77916"],
5     city="Montreal"
6 )
7
8 df_venues = df_venues.union(spark.createDataFrame([new_venue]))
```

4.1.4. Insert a new article in his author dataframe

Through this command we inserted the article created in section 4.1.2 to its authors. In order to do that, we selected the authors through the ids and we add the article id to their field `articles`.

Note: one of the author is the one created in section 4.1.1.

```
1 df_authors = df_authors.withColumn(
2     "articles",
3     f.when(f.col("_id") == "638db170ae9ea0d19fad7a79",
4         f.array_union(df_authors.articles, f.array(f.lit("
5             638db237d794b76f45c77916"))))\
6     .when(f.col("_id") == "638db170ae9ea0d19fad7a7a",
7         f.array_union(df_authors.articles, f.array(f.lit("
8             638db237d794b76f45c77916"))))
9     .otherwise(f.col("articles"))
10 )
```

4.1.5. Update the number of citations of referenced publications

Through the following snippet of code is possible to increment the `n_citations` field of the *Publications* referenced by the article created in section 4.1.2.

Note: field `n_citations` is updated for both the referenced articles.

```

1 df_articles = df_articles.withColumn(
2     "n_citation",
3     f.when(f.col("_id") == "53e99fe4b7602d97028bf743",
4         df_articles.n_citation+1) \
5     .when(f.col("_id") == "53e99fddb7602d97028bc085",
6         df_articles.n_citation+1)
7     .otherwise(f.col("n_citation"))
8 )

```

4.2. Queries

We have identified the following queries in order to show the system's basic functionalities. In the following sections title we wrote the basic requirements for every query, that, for ease of read, are represented as SQL clauses.

4.2.1. Query 1 - WHERE, JOIN

This query returns the type of the venue of an article with a the following title: *"Locality Sensitive Outlier Detection: A ranking driven approach"*.

```

1 df_articles.join(df_venues, df_articles.venue_raw == df_venues.raw, "
2     inner")\
3     .filter(f.col("title") == "Locality Sensitive Outlier
4         Detection: A ranking driven approach
5         ").select("title", "raw", "type").
6         show()

```

4.2.2. Query 2 - WHERE, LIMIT, LIKE

This query returns the articles whose title string contains *"Machine Learning"*.

```

1 df_articles.filter(f.col("title").like("%Machine Learning%")).limit(3).
2     show()

```

4.2.3. Query 3 - WHERE, IN, NESTED QUERY

This query finds authors that has the same nationality of at least one of the authors of *"Locality Sensitive Outlier Detection: A ranking driven approach"* article.

Description: we started creating the list of nationalities of the article's authors adding them to a list: `nationalities_list`. Then we look for authors with the a nationality contained in the list.

```

1 nationalities_list = df_articles.filter(f.col("title") == "Locality
    Sensitive Outlier Detection: A
    ranking driven approach")\
2     .select(f.explode(df_articles.authors.
    idAuth).alias("idAuth"))\
3     .join(df_authors, on=f.col("idAuth") ==
    df_authors._id)\
4     .select("nationality")\
5     .agg(f.collect_set("nationality")).
    collect()[0][0]
6
7 df_authors.filter(f.col("nationality")\
8     .isin(nationalities_list)).limit(10).show()

```

4.2.4. Query 4 - GROUP BY, JOIN, AS

This query finds the 3 most frequent keywords of articles written by italian authors.

Description: we started from finding all the italian authors, then we grouped them by articles. Then we count the keywords putting them in descending order.

```

1 df_italian = df_authors.filter(f.col("nationality") == "it")\
2     .select(f.explode("articles")).withColumnRenamed(
    "col", "articles")
3 df_italian = df_italian.groupby("articles").count()
4
5 df_keywords = df_italian.join(df_articles, df_italian.articles ==
    df_articles._id, "inner")\
6     .select("articles", f.explode("keywords")).
    withColumnRenamed("col", "keywords")\
7     .groupby("keywords")\
8     .agg(f.count("keywords").alias("n_occurences"))\
9     .sort("n_occurences", ascending=False)\
10    .limit(3).show()

```

4.2.5. Query 5 - WHERE, GROUP BY

This query finds the cities with more than 65 venues.

```
1 df_venues.groupby("city")\
2     .count()\
3     .filter(f.col("count") > 65)\
4     .sort("count", ascending=False).show()
```

4.2.6. Query 6 - GROUP BY, HAVING, AS

This query finds the field of studies that appears more than 15 times.

```
1 df_articles.select("_id", "title", f.explode("fos")).withColumnRenamed("
2                                     col", "fos")\
3     .groupby("fos")\
4     .agg(f.count("fos").alias("n_occurence"))\
5     .filter(f.col("n_occurence") > 15).show()
```

4.2.7. Query 7 - WHERE, GROUP BY, HAVING, AS

This query finds all the volumes with at least 5 articles in the dataset, published after 2000.

Description: we started from filtering the year of publication, then we grouped by `venue_raw` and `volume` in order to separate different edition of the same venue. After that we filtered the number of articles for every edition.

```
1 df_articles.filter(f.col("year") > 2000)\
2     .groupby("venue_raw", "volume")\
3     .agg(f.count("volume").alias("num_articles"))\
4     .filter(f.col("num_articles") > 4)\
5     .show()
```

4.2.8. Query 10 - WHERE, GROUP BY, HAVING, 2 JOINS

This query finds all the authors that published on more than 2 Journals.

```
1 df_exploded_authors = df_authors.alias("auth")\
```

```
2      .select("auth._id","auth.name", f.  
3          explode("auth.articles").alias("article"))\  
4      .join(df_articles.alias("art"), on=f.col("article") == df_articles._id)\  
5      .select("auth._id","auth.name","art._id",  
6          "art.venue_raw")\  
7      .join(df_venues.alias("ven"), on=f.col("venue_raw") == df_venues.raw)\  
8      .filter(f.col("type") == 1)\  
9      .groupBy("auth._id")\  
10     .agg(f.first("name").alias("name"),f.  
        countDistinct("raw").alias("venue_count"),f.concat_ws(" - ",f.  
        collect_set("raw")).alias("venues_list"))\  
        .filter(f.col("venue_count") > 2)\  
        .orderBy("venue_count", ascending=False)  
        .show(3,truncate=False)
```